

Text Mining Lessons Learned from Real Applications

ASIAS Technology and Tools Symposium
July 27, 2009

Elder Research, Inc.
300 W Main St., Suite 301
Charlottesville, Virginia 22903
434-973-7673
www.datamininglab.com
Elder@datamininglab.com



Outline

- Humans v. Machines?
- Mining tables and text for valuable information
 - Application areas at Elder Research
 - Ex: IRS Fraud detection
- Example **Government text mining** projects
 - Risk Profiling for NSA
 - Prioritizing CBP searches
 - Quick decisions for SSA disability
 - Document discovery for NGIC
 - Disease discovery for NCMI
- Factors for data/text mining project success



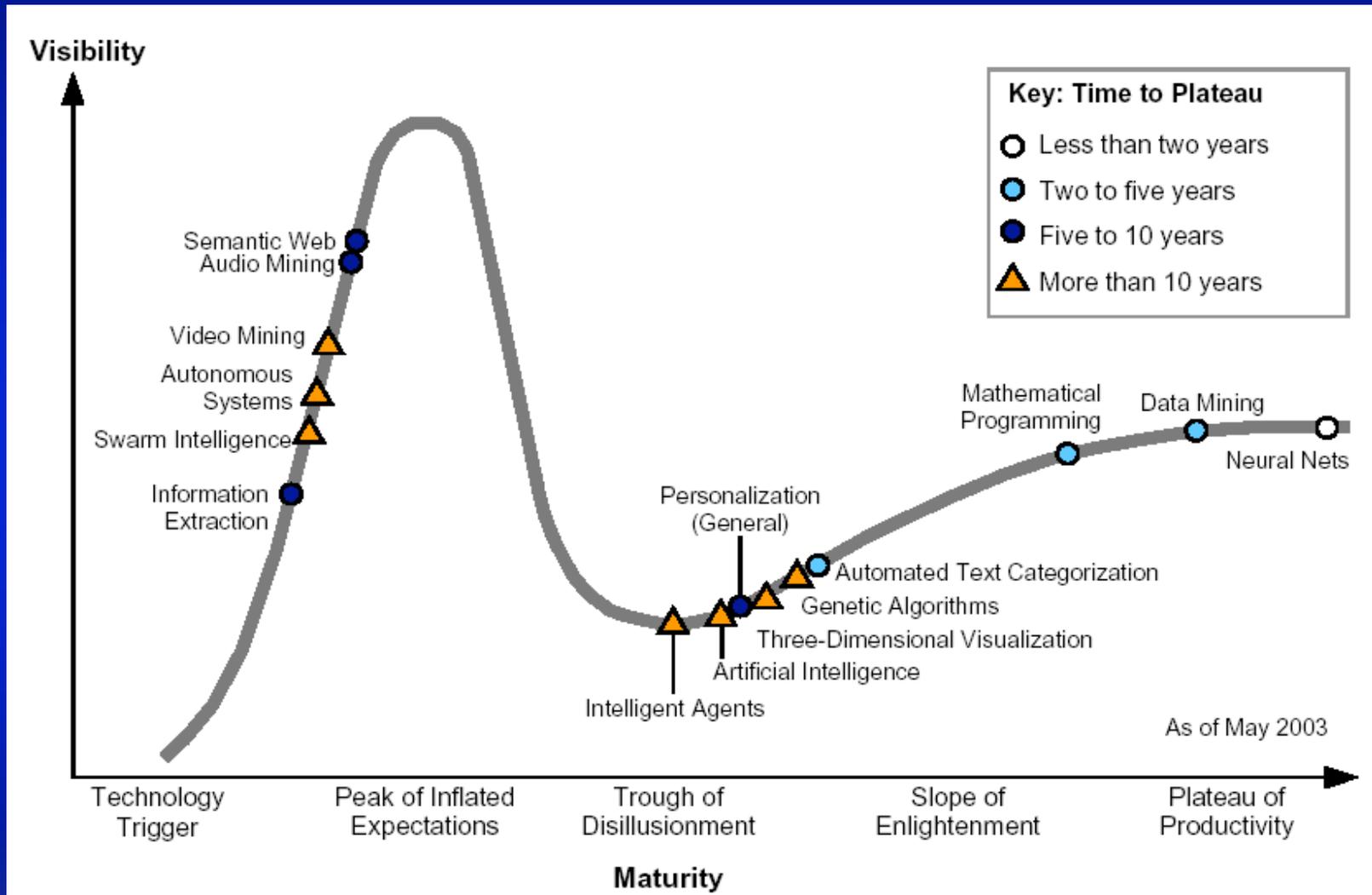
“Of course machines can think. After all, humans are just machines made of meat.”

- MIT CS professor

Our view:

Human and computer strengths are more complementary than alike.

Data Mining and the Hype Cycle



Selected ERI Clients

<p>Text Mining</p> <ul style="list-style-type: none"> • Social Security Administration • National Ctr. for Medical Intelligence • US Customs and Border Protection • National Ground Intelligence Center • Digital Reasoning Systems • US Dept. Education, NIDRR • Northrup Grumman • Ctr. for Navy SEAL Operations • US Army Tank Automotive RDEC 	<p>Investment</p> <ul style="list-style-type: none"> • Westwind Foundation • Oppenheimer Funds • CPTR • Vantage Consulting Group • Two Rivers Capital Mgmt. • N. de Rothschild Holdings • SAC Capital Advisors • R-Squared Trading, LLC • Energy Service Providers 	<p>Biometric & Pharmaceutical</p> <ul style="list-style-type: none"> • AstraZeneca • Pharmacia & UpJohn • SmithKline Beecham Pharm. • Rio Grande Medical Technology • Epsilon Group • Lumidigm • Dekalb Genetics • VeriLight
<p>CRM & Cross-Selling</p> <ul style="list-style-type: none"> • Hewlett Packard • AAA Michigan • Woodworker's Supply • Subscription Partners • HSBC 	<p>Behavioral & Web Analytics</p> <ul style="list-style-type: none"> • PDI Solutions • 3eDC • Richmond Police Dept. <p>Collaborative Filtering</p> <ul style="list-style-type: none"> • FindMoreFives.com (Netflix) 	<p>Tool Evaluations</p> <ul style="list-style-type: none"> • Defense Finance & Accounting Service • Fair Isaac • Electronic Warfare Assoc. <p>Reliability</p> <ul style="list-style-type: none"> • Lockheed Martin
<p>Entity Extraction & Link Analysis</p> <ul style="list-style-type: none"> • Department of Defense • Army Labs <p>Data Mining Technology</p> <ul style="list-style-type: none"> • Lucent 	<p>Data Warehousing & Business Intelligence</p> <ul style="list-style-type: none"> • Darden Solutions • Albemarle County • Georgetown University <p>Image Recognition</p> <ul style="list-style-type: none"> • Anheuser-Busch 	<p>Credit Scoring</p> <ul style="list-style-type: none"> • Capital One • Dealer Services • Grupo Ficohsa (Honduras) • Dollar Financial Group
<p>Training & Courses</p> <ul style="list-style-type: none"> • Dozens of Corporations, Government Agencies, Non-Profits, Universities, and Professional Associations 	<p>Optimization & Simulation</p> <ul style="list-style-type: none"> • Peregrine Systems • Westwind Foundation • NuTech Solutions • Finch Asset Mgmt. (Bermuda) • Commonwealth Comp. Research 	<p>Fraud Detection</p> <ul style="list-style-type: none"> • Internal Revenue Service • Defense Finance & Accounting Service • Federal Data Corporation • Mantas • Hewlett Packard

Example: Government, Financial Enforcement,
Large-Scale Production System, Strong Success

ERI Support to IRS Office of Refund Crimes

- ERI served as technical lead on a multi-year program to improve the detection of refund fraud in individual tax returns
 - Objective was to prioritize returns based on likelihood of fraud in order to optimize scarce resources (i.e., investigative analysts)
- Data mining system far out-performed existing methods
 - Found orders of magnitude more fraud than previous approach
 - Reduced workload with increased reliability, allowing analysts to focus on most promising and challenging cases
 - Improved refund time for honest filers because fewer “good” returns were held for initial review
- Resulting system was rolled-out nationwide ahead of schedule and with no adverse impact on existing workflow processes

ERI Support to US Customs & Border Protection Targeting & Analysis Systems Program Office

- Detecting unusual patterns of activity in land border crossings to identify criminal activity
- Mining free-form text to accurately identify commodities in sea-going transport
 - Mapping free-form text descriptions to established CBP codes which are used to identify unusual patterns of behavior

Social Security Administration

- President inquired why the SSA can't more quickly approve applicants with severe diseases
- Process can take up to 2 years
- Applicants need to be poor and unable to work (anywhere)
- Half of appealed applications are overturned
 - 5 layers of appeal
 - applicants actively worsening
- Commissioner Barnhart promised that SSA would use 21st century technology for Quick Approval of a sure subset

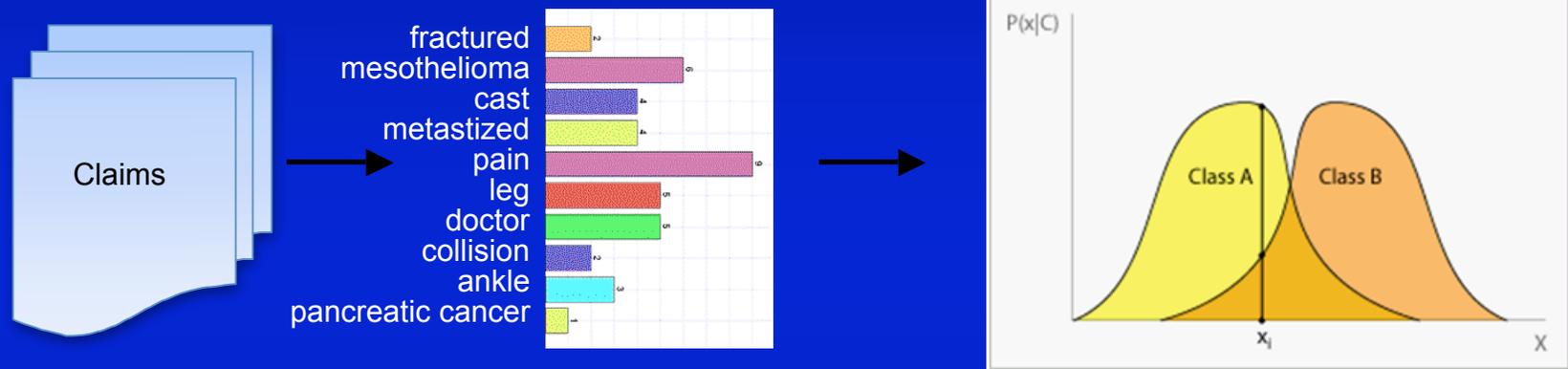
Text Mining Challenges

- Multi-Word Phrases (Concepts or Lemmas)
- Stemming (ex: Learn = Learning, Learned, Learns...)
- Synonyms (ex: ALS = Lou Gherig's Disease)
- Misspellings

Ex: 51 phrases found (by SPSS Clementine) for “Learning Disability”

learning disablitiy, learning deisability, learning disability, learining disabilities, learning disabiality, learningdisabilty, learning disabilty, learning disability, learning disabily, learning disabety, learning disability, learnoing disability, learning disabilities, learning disabiltiy, learning dsblty, learning disability, learnings disabilty, learningdisability, larning disabilities, learning disabilities, learning disabilities, learning disabilitties, learning disabilities, learning diasability, learning dasability, learnning disability, learning disabilities, lerning disability, learning disabilites, learneing disability, learninig disability, learning disaibilities, lernaning disability, learning disaibility, learnings disability, learning disabilitys, learning disabillity, learnings disabilities, learning diasability, learning disabiliites, learning dsiability, learning disabliity, learning disibilty, learning disabilities, learning disbality, learning disbility, learning disabilit, learningdisabilities, learningi disability, lerniung disabilities, learning disabliities, learning disaability, learning disabilities

Technology Approach at SSA: Bag of Words / Density-based Classification



$P(C|X)$ is the probability of case X being in class C .

Assign the class -- considering P and the *costs of misclassifications* -- that leads to the *minimum cost* decision.

30% baseline -> 90% model accuracy

Inductive Modeling Approach: Label cases for training (using analysts)

The screenshot shows the Altova XMLSpy interface. The main window displays XML markup for a text document. The text is segmented into three snippets, each containing collocation terms. The first snippet is about Baltimore's history, the second is about Birkhoff's mathematical work, and the third is about public schools. The XML tags used are `<Snippet>`, `<Collocation>`, and `<Noise>`. The interface includes a toolbar with options like 'Text', 'Schema/WSDL', 'Authentic', and 'Browser'. The status bar at the bottom indicates 'Ln 13, Col 900' and 'CAP NUM SCRL'.

Altova XMLSpy - [Collocation Example.xml *]

apr . 1861 was attacked by a mob . a disastrous fire in 1904 destroyed almost the entire downtown but enabled the emergence of a better planned city. In world wars i and ii baltimore was an important shipbuilding and supply-shipping center. during the 1960s and 70s however baltimore decayed rapidly losing population and commerce largely to neighboring suburbs . `<Collocation>` urban redevelopment `</Collocation>` in the late 1970s and 1980s included the construction of `< Collocation>`harborplace shops `</Collocation>` and restaurants in the `<Collocation>`inner harbor area `</Collocation>` the `<Collocation>`national aquarium `</Collocation>` shopping pavilions hotels a `</Snippet>`

`<Snippet>``<Collocation>`to abstract mathematics `</Collocation>` the teaching of mathematics and `< Collocation>`mathematical physics `</Collocation>` . from 1934 on he developed the concept of a lattice a `<Collocation>`generalized algebra `</Collocation>` with two operators and showed how a number of subjects e.g. `<Collocation>`boolean algebra `</Collocation>` `<Collocation>`projective geometry `</Collocation>` and `<Collocation>`affine geometry `</Collocation>` could be treated as special types of lattices . his text a `<Collocation>`survey of modern algebra `</Collocation>` with `< Collocation>`saunders maclane `</Collocation>` 1941 became a standard undergraduate his `< Collocation>`lattice theory `</Collocation>` `<Noise>`1940 3d ed . 1967 .2 `</Noise>` `<Collocation>` birkhoff george david `</Collocation>` 18841944 american mathematician `<Collocation>`b. overisel mich `</Collocation>` . . father of `<Collocation>`garrett birkhoff `</Collocation>` . the son of a physician he was educated at harvard `</Snippet>`

`<Snippet>`the public schools violated the principle of equal protection under the law guaranteed by

Text Schema/WSDL Authentic Browser

CollocationMarkup.xsd Collocation Example.xml

XMLSpy v2006 sp2 U Registered to Zach Buckner (Elder Research) ©1998-2005 Altova GmbH Ln 13, Col 900 CAP NUM SCRL

Build Statistical Model from Labeled Data

Microsoft Excel - bigram.xls

File Edit View Insert Format Tools Data Window Help Adobe PDF Type a question for help

E528 at

	C	D	E	F	G	H	I	J	K	L	M	N	
	V1	V1 Co	V2	V2 Co	L	C(V1)	C(V2)	C(V1V2)	Mean O	STD Out	Mean In	STD In	
517	at	TRUE	shiat	FALSE		4628	22	1	8.66E+00	6.65E+01	1.69E+00	1.32E+00	sunni people kept firing AT SH
518	leave	FALSE	except	FALSE		515	1240	1	5.08E+00	9.00E+00	2.81E+00	1.16E+01	shall neither take nor LEAVE E
519	activity	FALSE	that	TRUE		109	25965	1	4.04E+00	8.17E+00	8.79E+00	5.69E+01	security plan for each ACTIVIT
520	al-masjid-ai-haram.	FALSE	.	TRUE		1	279751	1	1	0	8.73E+00	1.71E+02	any other mosque excepting A
521	_o_8	FALSE	.	TRUE	2	4	279751	4	4	0	8.73E+00	1.71E+02	l_c_o_8_l_q_o_8_..._l_l_
522	30-31	FALSE	.	TRUE	2	4	279751	2	1.33E+00	5.77E-01	8.73E+00	1.71E+02	. koran 13 . 30-31 . they also say
523	scatter	FALSE	you	TRUE		9	15384	1	1	0	1.51E+01	8.08E+01	other ways which will SCATTE
524	the	TRUE	corpus	FALSE		130279	3	1	1.25E+01	1.06E+02		1	. you may keep THE CORPUS
525	.	TRUE	3.36	FALSE		279751	1	1	2.07E+01	3.66E+02		1	from the outcast satan . 3.36 -
526	shawahid	FALSE	al-tanzil	FALSE	1	4	2	2	1.33E+00	5.77E-01		2	v3 . p3716 SHAWAHID AL-TA
527	smelled	FALSE	him	TRUE		2	22843	1	1	0	1.81E+01	2.70E+02	and kissed him and SMELLEC
528	matters	FALSE	at	TRUE		143	4628	1	2.88E+00	4.99E+00	3.36E+00	1.34E+01	and i discussed such MATTEI
529	those	FALSE	surrendered	FALSE		4254	47	1	9.27E+00	1.12E+02	2.24E+00	4.44E+00	but one house of THOSE SUF
530	other	TRUE	soul	FALSE		2762	331	1	3.47E+00	1.53E+01	5.61E+00	1.20E+01	at all for any OTHER SOUL . th
531	real	FALSE	prophecy	FALSE		147	28	1	1.59E+00	1.26E+00	3.11E+00	3.95E+00	real interpretation of this REA
532	not	FALSE	masterminded	FALSE		14173	1	1	7.14E+00	3.06E+01		1	the attacks applauded if NOT I
533	of	TRUE	architecture	FALSE		76432	38	8	1.11E+01	2.29E+02	1.81E+00	1.81E+00	issue of the use OF ARCHITE
534	imam	FALSE	in	TRUE		1328	31692	15	7.78E+00	3.13E+01	5.44E+00	4.54E+01	the qur'an by the IMAM IN the
535	.	TRUE	thenceforth	FALSE		279751	2	1	2.07E+01	3.66E+02		1	0 bright gleam of dawn . THENC
536	ai-isra'	FALSE	.	TRUE		1	279751	1	1	0	8.73E+00	1.71E+02	an-nisa' . 163 & AI-ISRA' . 55 pe
537	weak	FALSE	commodity	FALSE		171	25	1	2.80E+00	6.83E+00	1.67E+00	1.50E+00	rains in 2001 . WEAK COMMO
538	was	TRUE	listening	FALSE		13895	58	6	6.82E+00	3.40E+01		2	1.69E+00 a man and imran 'WAS LISTEN
539	my	TRUE	auliya'	FALSE		4914	15	2	5.55E+00	2.07E+01	1.50E+00		7.07E-01 helps and protectors . MY A
540	us	FALSE	operations	FALSE		3579	30	1	7.41E+00	4.80E+01	1.30E+00	6.57E-01	countries to condemn the US
541	benhamou	FALSE	.	TRUE		1	279751	1	1	0	8.73E+00	1.71E+02	initiatives for development mc
542	when	TRUE	protecting	FALSE		7404	76	1	1.33E+01	9.25E+01	2.42E+00	2.73E+00	person during self-defense or
543	labors	FALSE	among	FALSE		4	1872	1	1.33E+00	5.77E-01	2.96E+00	9.86E+00	labor of any that LABORS AV
544	muslims	FALSE	alone	FALSE		1285	429	1	4.84E+00	2.20E+01	3.30E+00	8.41E+00	not belong to the MUSLIMS A
545	considered	FALSE	praiseworthy	FALSE		222	15	1	2.31E+00	4.60E+00	1.88E+00	2.10E+00	possessors of which are CON
546	fis's	FALSE	armed	FALSE		1	65	1	1	0	1.91E+00	2.42E+00	by the late-1990s and FIS'S AF

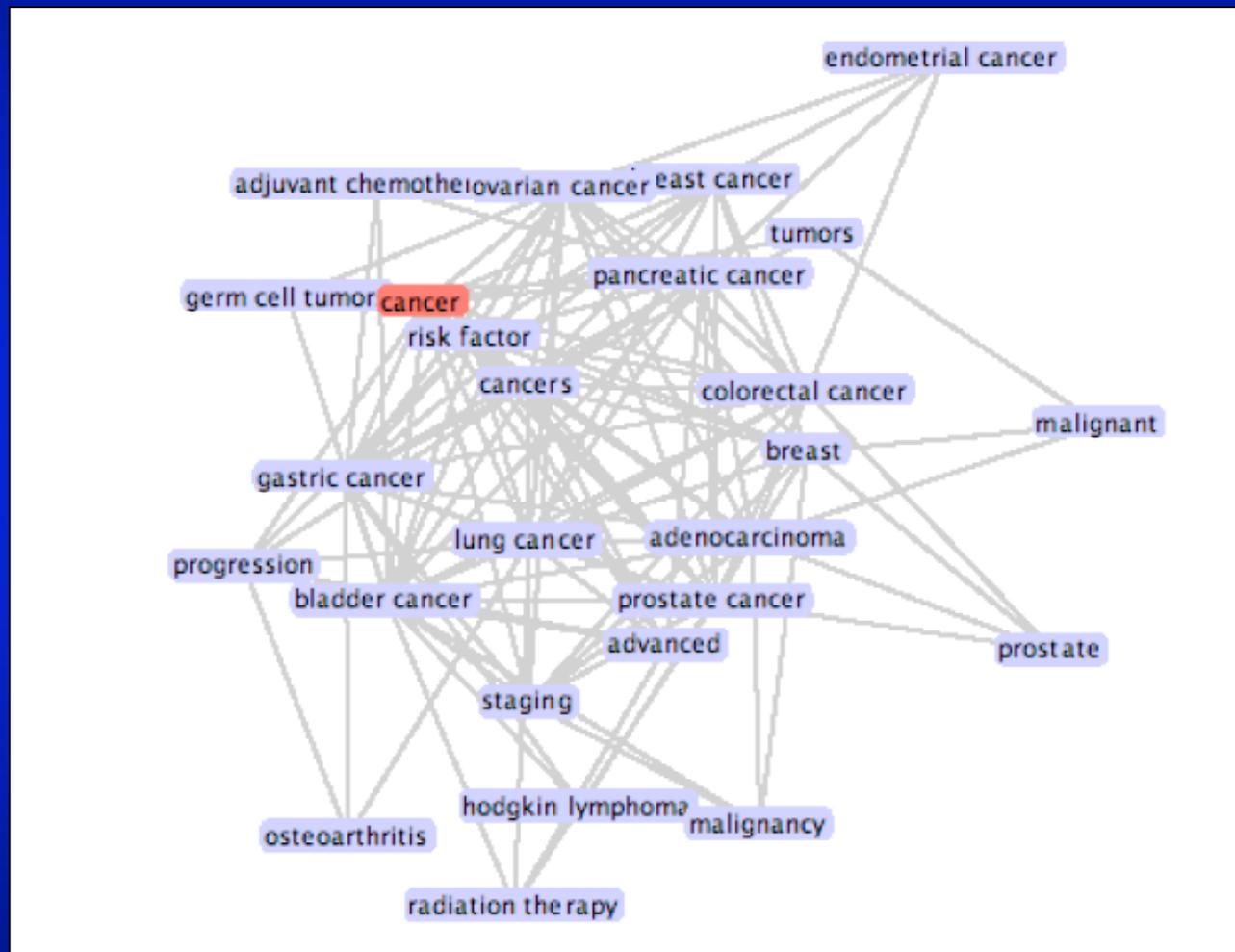
Bigrams Common Training

Ready

Apply Model to New Data

((Ahmed Omar) (Saeed Sheikh)) Known as "**Sheikh Omar**", ((Ahmed Omar) (Saeed Sheikh)) is a (British citizen) of Pakistani descent, with links to various Islamic-based (terrorist organisations), including Al-Qaeda and Harkat-ul-Mujahideen. He has been mentioned in many **conspiracy theories** linking him with the CIA and the ISI, Pakistan's **intelligence agency**, for whom he was allegedly an informer. Many of the sensationalist and inaccurate reports - the basis for the **conspiracy theories** involving him - arose in the confusion of the early weeks after the (9/11 attacks), when investigators did not have a (clear picture) of the plot, and followed a great many leads that later turned out to be (dead ends). As of 2005 the suspicion that Sheikh (played a part) in the funding of the (9/11 attacks) has been fading. However, as the factual support for his involvement has crumbled, **conspiracy theories** about him have thrived. Much of what follows is perhaps best seen in that context. In his youth he attended **Forest School Snaresbrook**, a (public school) in (North-East London), whose alumni include England cricket captain **Nasser Hussain**. He also attended the prestigious **London School of Economics**. (The Times) describes **Saeed Sheikh** as "no ordinary terrorist but a man who has connections that reach high into Pakistan's (military and (intelligence) elite) and into the innermost circles of **Osama Bin Laden** and the (al-Qaeda organization)." According to ABC, Sheikh began working for the ISI in 1993. By 1994 he was operating (terrorist training camps) in Afghanistan and had earned the title of (bin Laden's) "special son." At the time, the Taliban were beginning to dominate Afghanistan, much due to support received from the ISI. In (May 2002), the **Washington Post** quotes an unnamed Pakistani as saying that the ISI paid Sheikh's legal fees during his 1994 trial in India on charges of kidnap. However, this claim has not been confirmed by any other source. An unnamed senior-level **U.S. government** (source told) CNN in (October of 2001) that (U.S. investigators had discovered) that someone using the alias (Mustafa (Muhammad Ahmad)) possibly ((Ahmed Omar) (Saeed Sheikh)), allegedly a long-time ISI informer, had sent about \$100,000 from the **United Arab Emirates** to **Mohammed Atta**, the suspected hijack ringleader of the (September 11, 2001 attacks). "Investigators said Atta then distributed the funds to conspirators in Florida in the weeks before the deadliest (acts of terrorism) on U.S. soil that destroyed the **World Trade Center**, (heavily damaged) the Pentagon and (left thousands) dead. In

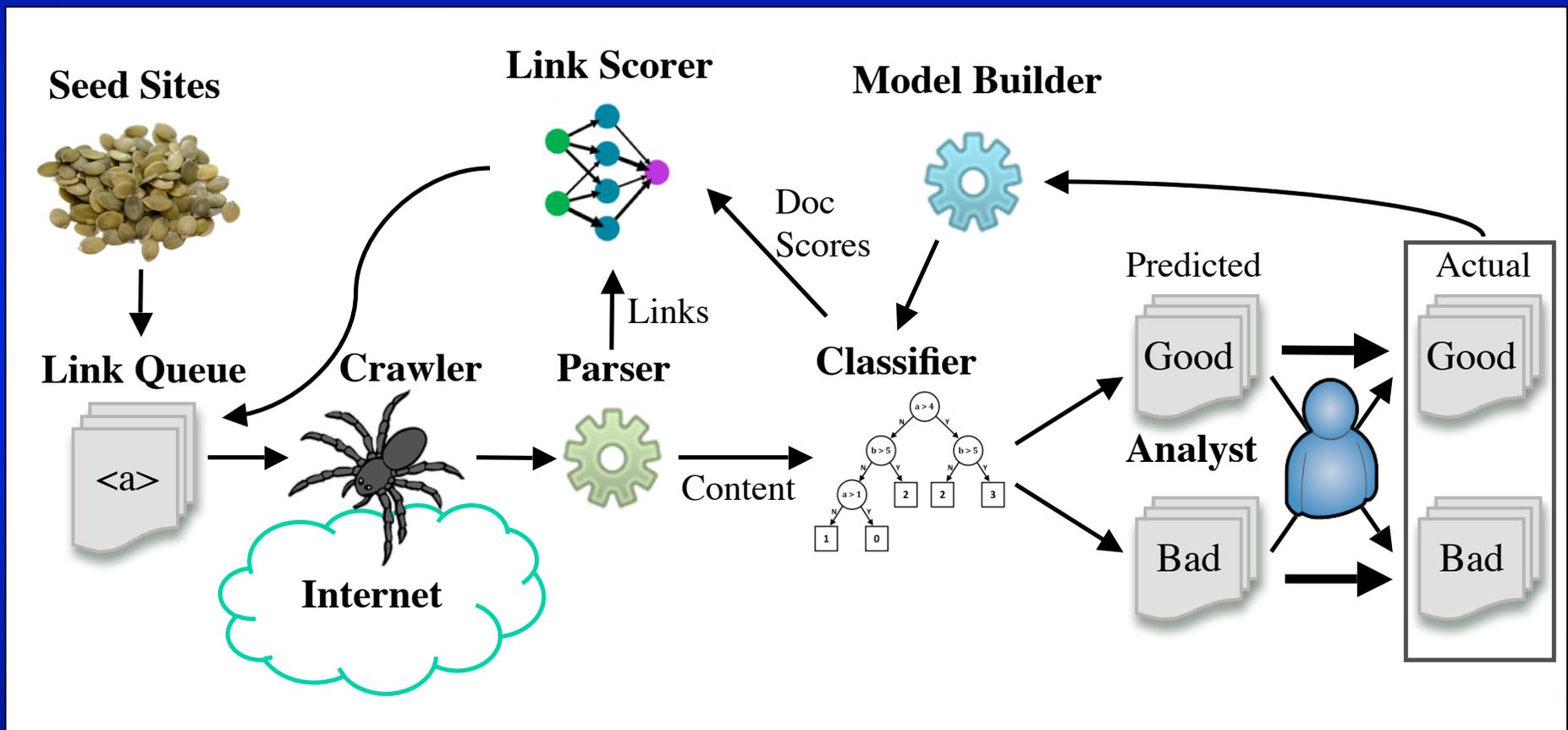
Collocation & Association Networks (National Ground Intelligence Center)



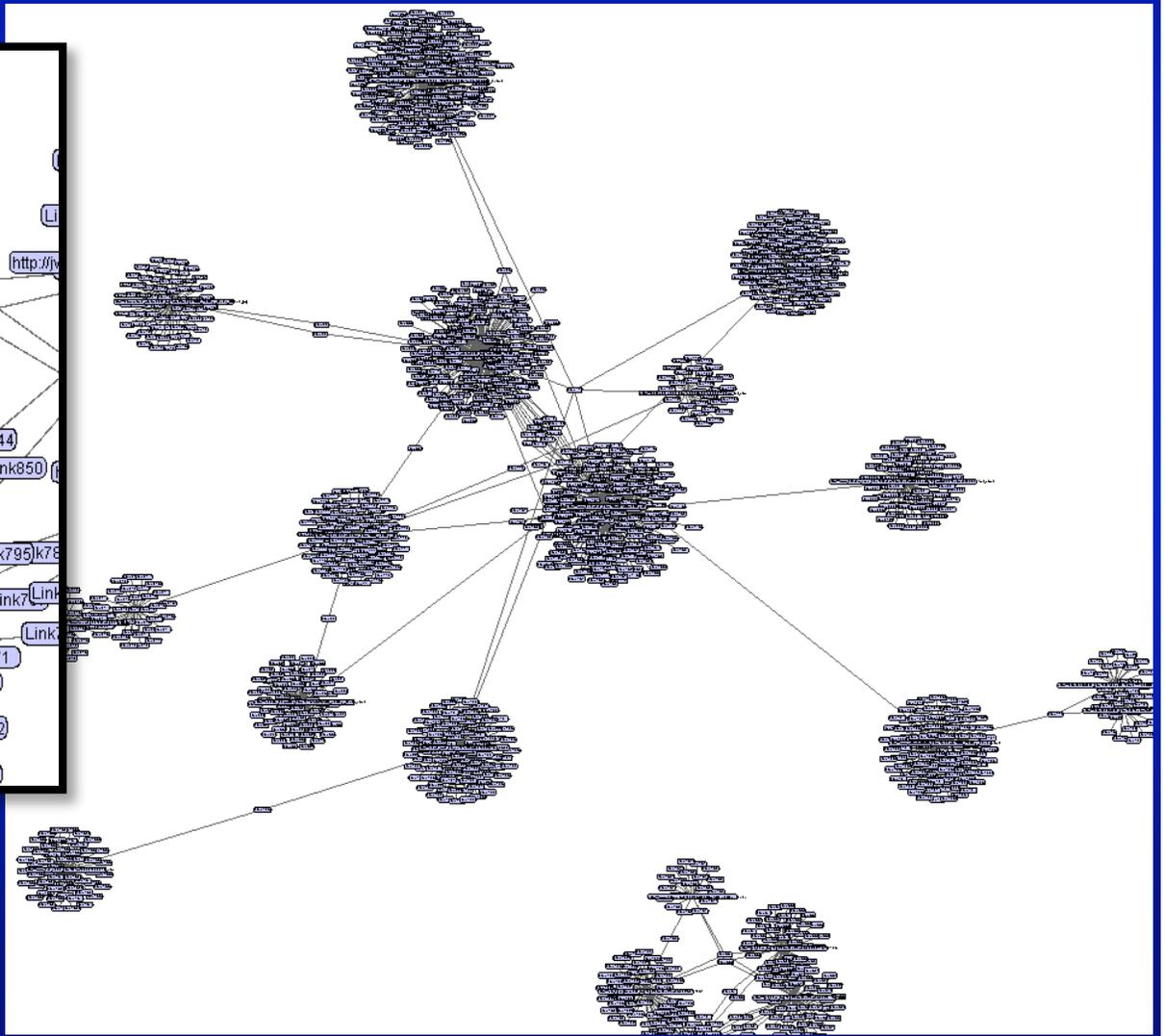
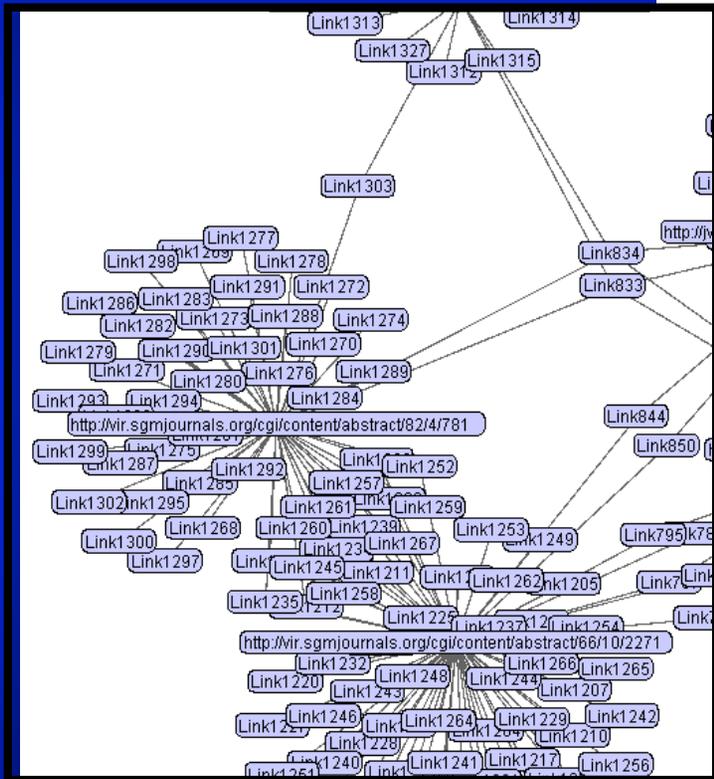
Improving Search with Collocation and Association

The screenshot displays the Yahoo! search page. At the top left is the Yahoo! logo. The search bar contains the text "cpu mac mini" and a "Web Search" button. Below the search bar, a dropdown menu titled "Also search for:" lists several suggestions: "macbook pro", "macbook", "imac", "ibook", "fedora core", "quantum computer", "macintosh", "mbp", "toolkit", and "hardware". The "macintosh" suggestion is highlighted with a mouse cursor. To the left of the search bar are links for "My Yahoo!" and "My Mail". Below the search bar are navigation tabs for "Web", "Images", "Video", "Local", "Shopping", and "more". On the left side of the page, there is a vertical menu with icons and labels for "Answers", "Autos", "Finance", "Games", "GeoCities", "Groups", "HotJobs", "Maps", "Movies", "Music", "Personals", "Real Estate", "Shopping", and "Sports". The main content area features a "Featured" section with a date of "Jul 11, 2007" and a headline "Heart-h...". Below the headline is an image of pistachio nuts and a sub-headline "Eating pistachio nuts lowers your blood pressure". To the right of the featured article are several smaller news items with thumbnails and headlines, such as "Double-amputee to run in foot race" and "Is new 'Potter' film his darkest adventure yet?". On the right side of the page, there is a personalized greeting "Hi, John" with a "Sign Out" link. Below the greeting are several utility buttons: "Mail", "Messenger", "Radio", "Weather 90°F", "Local", and "Horoscopes". At the bottom of the page, there is a section for "In the News" with tabs for "World", "Local", and "Video".

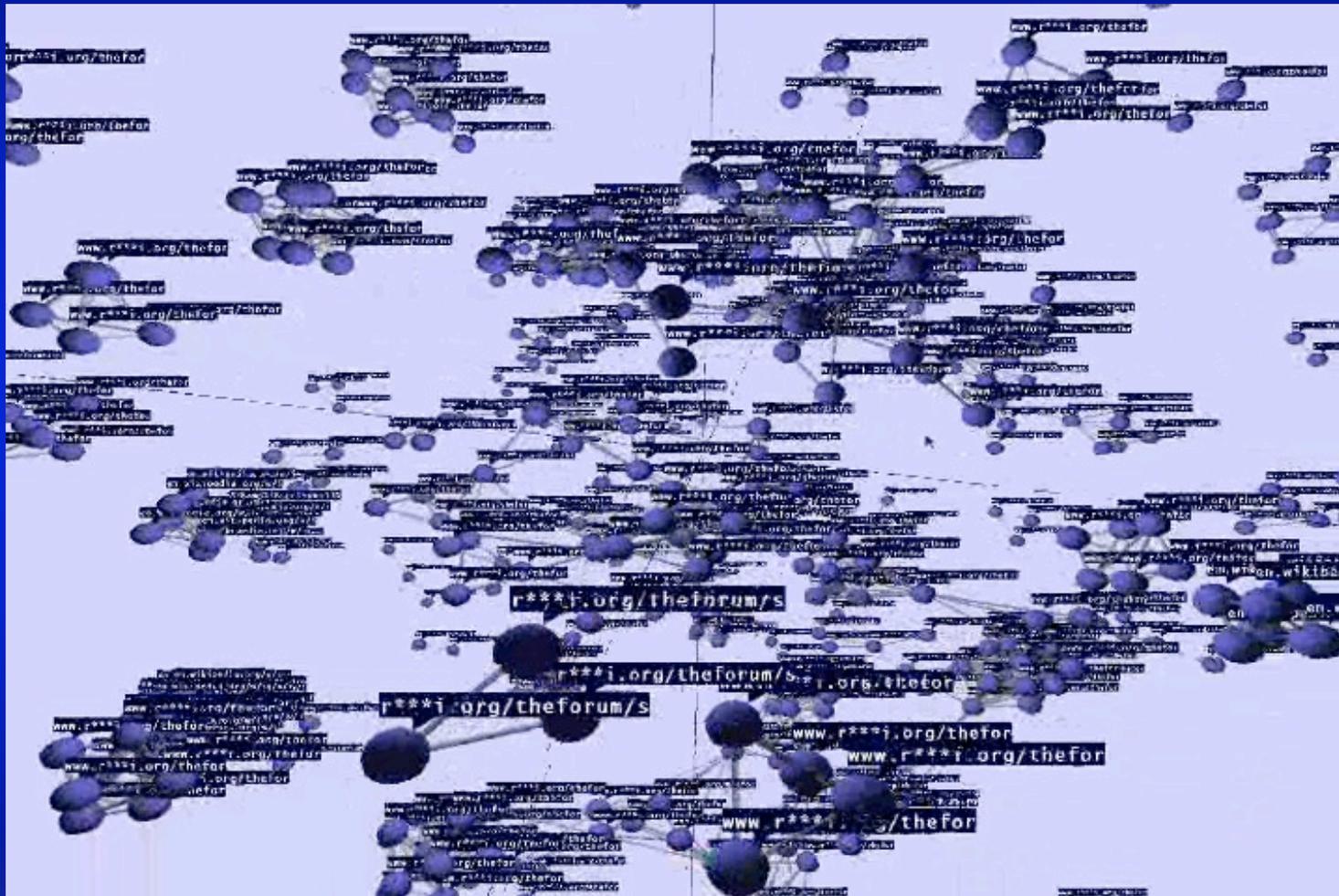
Monitoring the Web for Infectious Animal Diseases (National Center for Medical Intelligence)



Exploring Web of References



Visualize Document Similarity (ERI Tool)



Prioritize Documents for Analyst Attention through Discovered Clusters



Necessary Ingredients for Data/Text Mining Project Success

- Gain Expected: either:
 - Leverageable - an incremental improvement will matter, or
 - “Low-hanging fruit” - nobody’s yet dared attack the problem
- Interdisciplinary Team: experts needed in business area, statistics, algorithms, and databases
- Data Vigilance: capture and maintain the accumulating information stream
- Time: learning occurs over multiple cycles
- Business Champion is essential

How to Manage Data/Text Mining Projects

- Assess data assets (what treasure could be hidden in our sludge?)
 - Data caretaker must be on board
- Identify pain points in current production process
 - What improvements would have the biggest impact?
- Brainstorm ideal process
 - External expertise extremely efficient here
- Conduct a pilot project. Simultaneously:
 - “Hit a single”: automate key task, create dashboard, or graphic
 - “Swing for fences”: attack core weakness
- Have key staff work closely with analytic experts
 - Transfer technology
 - Internalize essential steps
- Prove ROI. Make allies and decision-makers look good

John F. Elder IV

Chief Scientist, Elder Research, Inc.



DR. JOHN ELDER HEADS A DATA MINING CONSULTING TEAM WITH OFFICES IN CHARLOTTESVILLE, VIRGINIA AND WASHINGTON DC (WWW.DATAMININGLAB.COM). FOUNDED IN 1995, ELDER RESEARCH, INC. FOCUSES ON INVESTMENT, COMMERCIAL AND SECURITY APPLICATIONS OF ADVANCED ANALYTICS, INCLUDING TEXT MINING, STOCK SELECTION, IMAGE RECOGNITION, PROCESS OPTIMIZATION, CROSS-SELLING, BIOMETRICS, DRUG EFFICACY, CREDIT SCORING, MARKET TIMING, AND FRAUD DETECTION.

JOHN OBTAINED A BS AND MEE IN ELECTRICAL ENGINEERING FROM RICE UNIVERSITY, AND A PHD IN SYSTEMS ENGINEERING FROM THE UNIVERSITY OF VIRGINIA, WHERE HE'S AN ADJUNCT PROFESSOR TEACHING OPTIMIZATION OR DATA MINING. PRIOR TO 14 YEARS AT ERI, HE SPENT 5 YEARS IN AEROSPACE DEFENSE CONSULTING, 4 HEADING RESEARCH AT AN INVESTMENT MANAGEMENT FIRM, AND 2 IN RICE'S *COMPUTATIONAL & APPLIED MATHEMATICS* DEPARTMENT.

DR. ELDER HAS AUTHORED INNOVATIVE DATA MINING TOOLS, IS A FREQUENT KEYNOTE SPEAKER, AND WAS CO-CHAIR OF THE 2009 *KNOWLEDGE DISCOVERY AND DATA MINING* CONFERENCE, IN PARIS. JOHN'S COURSES ON ANALYSIS TECHNIQUES – TAUGHT AT DOZENS OF UNIVERSITIES, COMPANIES, AND GOVERNMENT LABS – ARE NOTED FOR THEIR CLARITY AND EFFECTIVENESS. DR. ELDER WAS HONORED TO SERVE FOR 5 YEARS ON A PANEL APPOINTED BY THE PRESIDENT TO GUIDE TECHNOLOGY FOR NATIONAL SECURITY. HIS BOOK ON PRACTICAL DATA MINING, WITH BOB NISBET AND GARY MINER, APPEARED IN MAY 2009.



JOHN IS A FOLLOWER OF CHRIST AND THE PROUD FATHER OF 5.